

Shotgun Metagenomic and Physicochemical Characterisation of Soil Microbial Communities from the UNISEL Bird Sanctuary, Malaysia

Hasdianty Abdullah^{1,2*}, Mohd Fadzli Ahmad^{1,2}, Fridelina Sjahrir¹,
Moohamad Ropaning Sulong³, Hazeeq Hazwan Azman¹,
Maegala Nallapan Maniyam⁴, and Nor Suhaila Yaacob^{1,2}

¹Faculty of Engineering and Life Sciences, Universiti Selangor, Jalan Timur Tambahan, 45600 Bestari Jaya, Selangor, Malaysia

²Institute of Bio-IT Selangor, Universiti Selangor, Jalan Zirkon A7/A, Seksyen 7, 40000 Shah Alam, Selangor, Malaysia

³Halalan Thayyiban Research Centre (HTRC), Universiti Islam Sultan Sharif Ali (UNISSA), Kampus Sinaut KM 33, Jalan Tutong, Kampung Sinaut, 1741 Tutong TB, Brunei Darussalam

⁴Faculty of Health Sciences, Universiti Selangor, Jalan Zirkon A7/A, Seksyen 7, 40000 Shah Alam, Selangor, Malaysia

ABSTRACT

This study explores the microbial diversity within the newly discovered UNISEL Bird Sanctuary using a metagenomics approach. The objective was to establish the first baseline shotgun metagenomic profile of the UNISEL Bird Sanctuary soil microbiome to support future ecological, environmental, and conservation research. Soil samples were subjected to whole-genome shotgun sequencing to characterise microbial community structure and functional potential. Approximately 2.85 million contigs were generated, predicting 3,282,381 genes. Taxonomic profiling using Kraken 2 showed that bacteria represented a major proportion of the community, accounting for approximately 30.01% of clean reads and 34.83% of assembled contigs. A large number of reads remained unclassified, indicating the presence of potentially novel microbial diversity. Functional

annotation successfully classified 44% of coding genes against databases including RefSeq, Swiss-Prot, InterProScan, eggNOG, MEGARes, and KEGG. Comparative metagenomic analysis with publicly available datasets revealed similarities and differences in taxonomic composition, Clusters of Orthologous Groups (COG) classification, and metal resistance genes. The UNISEL Bird Sanctuary soil sample (S3) exhibited higher abundances of copper, iron, and mercury resistance genes, suggesting adaptation

ARTICLE INFO

Article history:

Received: 05 May 2025

Accepted: 04 May 2026

Published: 29 May 2026

DOI: <https://doi.org/10.47836/pitas.49.3.02>

E-mail addresses:

dianty@unisel.edu.my (Hasdianty Abdullah)

fadzli@unisel.edu.my (Mohd Fadzli Ahmad)

fridelina@unisel.edu.my (Fridelina Sjahrir)

ropaning.sulong@unissa.edu.bn (Moohamad Ropaning Sulong)

hazeeq87@unisel.edu.my (Hazeeq Hazwan Azman)

maegala@unisel.edu.my (Maegala Nallapan Maniyam)

shuhaila@unisel.edu.my (Nor Suhaila Yaacob)

* Corresponding author

to specific environmental conditions. This study highlights the importance of documenting microbial diversity for ecosystem health monitoring, ecological understanding, and conservation planning. Further studies involving deeper sequencing and advanced analyses are required to better understand the functions and adaptive mechanisms of uncharacterised microbial communities in this ecosystem.

Keywords: Bioinformatics, metagenomics, microbial diversity, landfills, Whole-Genome Shotgun Sequencing (WGS)

INTRODUCTION

The bird sanctuary has become an increasingly attractive recreational area for tourists and environmentalists to engage in bird-watching activities. This piece of land serves as a home to migratory birds, in which birds are protected and encouraged to breed, as well as promotes the survival and rehabilitation of the birds. In addition, the bird sanctuary area offers natural facilities that serve the conservation of various species and their natural habitat. UNISEL Bird Sanctuary or *Tasik Lindungan Burung* UNISEL is a freshwater-related conservation site found within the Universiti Selangor campus of Selangor, Malaysia. The sanctuary has a mosaic of rehabilitated lakes and other surrounding land habitats created as a result of historical tin mining activities that have evolved through natural ecological succession. These environments support a variety of vegetation and offer significant habitat and stopover for both resident and migratory bird species. The UNISEL Bird Sanctuary is a unique tropical ecosystem to be studied in ecological and environmental studies due to its limited anthropogenic disturbance during recent years. The geochemical legacy of its past mining activity helps to create unique physicochemical characteristics of the soil that affect the dynamics of microorganisms and plant communities. Thus, the sanctuary is an important natural laboratory to investigate soil health, microbial diversity, ecosystem recovery, and preservation of biodiversity in freshwater post-mining landscapes.

Although they are ecologically important, no prior baseline data have been provided on the microbial community structure or functional gene composition of soil at the UNISEL Bird Sanctuary. This information gap restricts the capacity to evaluate the health of the ecosystems, identify the initial signs of environmental stress or pollution, and inform evidence-based conservation and management approaches. Thus, the detailed preliminary study combines physicochemical characterisation and shotgun metagenomic analysis to understand the microbial community structure, functional possibilities, and adaptive characteristics of soil microorganisms in this sanctuary.

Microorganisms play a significant role in mineral recycling and nutrient absorption by flora and fauna in the natural habitat and facilitate many activities, such as the recycling of essential elements and nutrients, biogeochemical cycles, and the development of soil structure (Fulke, Mahajan and Gothawal, 2026). Note that microorganisms are vital in

waste processing, growth-promoting, and the reproduction of plants and animals (Husain, 2022). Hence, exploring microbial communities' profiles and the existing diversity is an interesting scope. For this purpose, metagenomics will be a powerful culture-independent method for studying the microbes and for formulating a microbial diversity profile of uncultured microorganism access as compared to other microorganism profiling methods (Kumar, 2026).

This paper presents our current study of microbial diversity in UNISEL's newly discovered bird sanctuary through a whole-genome sequencing metagenomics approach. Documenting microbial diversity offers numerous benefits, particularly within the context of a bird sanctuary soil metagenomics study. In addition, it provides a baseline for monitoring ecosystem health by identifying the types and abundance of microbes present in the soil, enabling future comparisons to assess the impact of environmental changes or pollution. It also aids in understanding ecosystem functions, as microbial communities play crucial roles in nutrient cycling and decomposition processes. Detailed knowledge of microbial diversity also informs conservation and management strategies (Crawford et al., 2026; Li et al., 2025). As an example, the discovery of unfavourable microbes depletion or an increase in unfavourable microbes enables specific measures to be taken to achieve the restoration of ecological balance. The comparative metagenomics of the study, which compared the soil sample to other datasets, can assist in understanding the differences in the microbial community studied in comparison to other studies, which helps customise the conservation. Also, it can be used in the assessment of the environmental impacts as observed in the research on the copper, iron, and mercury metal resistance gene identification, which suggested potential adaptation to certain environmental factors. Therefore, the abundance and diversity of these genes can be monitored to give information on the amount and effects of the pollutants in the soil.

MATERIALS AND METHODS

Soil Sampling

Sampling was done at three spatially different locations in the UNISEL Bird Sanctuary, Selangor, Malaysia (3°25'52.1" N, 101°26'11.1" E) at a depth of 0 - 20 cm. Sampling sites were identified in undisturbed areas of the sanctuary to capture the indigenous soil microbial communities with minimal anthropogenic impact. Minimal anthropogenic disturbance was the main site selection criteria because the primary goal of this study was to record a baseline soil microbial diversity as a natural heritage and long-term ecological reference. The subsamples obtained were pooled and mixed thoroughly to create one composite soil sample that could represent the soil environment in the sanctuary. The samples of the soils were moved into sterile sampling containers, then moved to the laboratory in insulated ice boxes with ice packs and kept in the freezer at -20 °C until further examination.

Composite sampling has been widely embraced in exploratory soil metagenomics to lessen spatial heterogeneity and create representative baseline microbial profiles, especially in cases where shotgun metagenomic sequencing is constrained by cost. The same methods have been used in previous metagenomics research (Zainun et al., 2018; Xi et al., 2026; Salam & Obayori, 2026; McLaren et al., 2026), which validated its adaptability in measuring baseline biodiversity. Nevertheless, the existence of a single composite samples was used does not allow testing of the fine-scale spatial heterogeneity within the sanctuary; this will be resolved in the future by replicated and spatially stratified sampling designs.

DNA Extraction

A soil sample of UNISEL Bird Sanctuary was used to extract total Deoxyribonucleic Acid (DNA). Sterile tools were used to collect the soil to prevent contamination, and the sample was stored in sterile conditions at 4 °C until it was processed. The DNA extraction was done on the basis of the Qiagen PowerSoil DNA Isolation Kit as per the protocol of the manufacturer (Cat. No. 12888-100). Around 0.25 g of soil was put in a PowerBead Tube, and solution C1 was added to allow cell lysis. Mechanical lysis was performed on the sample by bead-beating to effectively lyse microbial cells. The resultant lysate was centrifuged, and the supernatant was transferred to a clean tube. The solutions C2 and C3 were added sequentially to eliminate inhibitory substances, such as humic acids, and centrifuged. Solution C4 and a silica membrane spin column were used to achieve DNA binding. Solution C5 was used to wash the column to make sure that the DNA is pure, and Solution C6 was used to elute the DNA. The extracted DNA was kept at -20 °C to be used downstream. After extraction, the DNA quantity was measured with Qubit® 3.0 Fluorometer (Invitrogen, Carlsbad, CA, USA) using the double-stranded DNA (dsDNA) High-Sensitivity (HS) assay. Moreover, quality was checked using agarose gel electrophoresis to confirm the integrity and lack of any major contaminants.

Shotgun Metagenomics Analysis

Shotgun metagenomics analysis is an effective method that is applied to study the taxonomic composition and functional capacity of microbial communities directly using environmental samples. It is also a very detailed study, which involves the sequencing of all the genetic material in a sample. This approach allows detecting a wide spectrum of microorganisms, such as bacteria, archaea, fungi, and viruses, and the characterisation of metabolic pathways and functional genes. Shotgun metagenomics was applied in this study to examine the structure of the microbial community, as well as the functional profiles of the extracted soil DNA.

Data QC and Pre-processing

SolexaQA++ and Bowtie2 were used to perform data quality control and pre-processing. First, the quality of raw sequencing reads was tested with SolexaQA++, which removed low-quality bases to retain high-quality sequences that could be analysed downstream. After that, Bowtie2 was run to align the reads with the PhiX genome, which is a frequent spike-in control in sequencing libraries to eliminate any PhiX contamination. Following the removal of PhiX sequences, paired-end reads were filtered and combined to produce high-quality forward and reverse sequences to be analysed later.

Metagenomics Assembly

Metagenomic assembly was carried out with MetaSPAdes, a dedicated assembler that is designed to assemble microbial genomes with metagenomic data. The resulting high-quality paired-end reads following pre-processing were fed into MetaSPAdes, which takes a multi-step strategy that involves a combination of error correction, graph-based assembly, and iterative refinements to create contiguous sequences (contigs). The assembly that was obtained offered a basis upon which taxonomic and functional analysis could be conducted.

Taxonomic classification and profiling

Kraken 2, a K-mer-based taxonomic classification tool, was used to perform taxonomic classification and profiling. The metagenomic assembly's high-quality reads were matched against the National Centre of Biotechnology Information (NCBI) non-redundant database, which contains extensive genomic reference sequences. Kraken 2 was used to assign taxonomic labels to each read, and this allowed the accurate identification and quantification of microbial taxa in the sample. In line with this, the output was handled to produce taxonomic profiles and relative abundance data, which provided detailed information on the microbial community structure.

Functional Annotation

Analysis of the metagenomic data was conducted through functional annotation of the data using various reference databases and tools to be comprehensive. The annotated sequences were matched to the databases of RefSeq and Swiss-Prot to identify high-confidence protein-coding genes. InterProScan was used to predict protein domains and functional motifs, and eggNOG was used in orthology-based functional assignments and Gene Ontology (GO) classification. Antimicrobial resistance genes were identified using the MEGARes database, which offers information on possible resistance mechanisms. To minimise false-positive assignments, a conservative E-value of 1×10^{-5} was used to cut off searches in BLAST-based searches against RefSeq, Swiss-Prot, and MEGARes.

The reconstruction of metabolic pathways and functional profiling was done through the Kyoto Encyclopedia of Genes and Genomes (KEGG) database, which helped in identifying the microbial metabolic capabilities. The multi-database methodology provided a precise and comprehensive functional characterisation of the microbial community. The use of this multi-database method was selected in particular to avoid the variable coverage rates of environmental metagenomes, which enable the reliability of determining the functional potential core of the soil despite the high number of novel or uncharacterised genes.

Comparative Analysis

Comparative analyses of the metagenomic data produced in this study (S3) with two publicly available datasets, SRR21870203 and SRR21870204, were done to assess the structure of the microbial community and to determine the functional profiles. Kraken 2 was used to taxonomically profile the datasets, and the results were compared to evaluate the diversity and composition of microbial communities. Clusters of Orthologous Groups (COG) analysis was used to provide functional categories that allowed comparison of genes. The frequency of metal resistance genes was compared by aligning annotated sequences of each dataset with the MEGARes database. The Principal Component Analysis (PCA) was conducted to analyse beta diversity to identify the differences between the microbial community structures across the datasets. Also, Canonical Correspondence Analysis (CCA) was applied to measure associations between taxonomic composition and soil physicochemical characteristics, as well as to measure associations between metal resistance genes and soil parameters. These multivariate analyses offered information on environmental and functional drivers of microbial community dynamics among the samples.

RESULTS AND DISCUSSION

Physicochemical Properties of Soil Samples

The main objective of this study was to examine the microbial diversity in the Bird Sanctuary soil of UNISEL (later named S3) using a metagenomics method. The physicochemical composition of the soil was first established to reveal the dynamic interaction between the properties of the soil and the microbial communities that thrived in it. Therefore, the knowledge of this interaction is essential to the evaluation of soil health, ecological role, and the possibility of agricultural or ecological use.

Table 1 provides a summary of the S3 sample physicochemical analysis, supplemented by the heavy metal limits established by the Department of Environment (DOE) Malaysia, which offers a detailed analysis of its chemical, physical, and biological characteristics. The parameters examined are macronutrients, trace metals, toxic elements, Total Organic Carbon (TOC), physical properties, and microbial richness obtained through metagenomic data.

Table 1
UNISEL bird sanctuary soil samples physicochemical results

No.	Parameter	Unit	Concentration	Standard Limit*
1.	Sodium, Na		0.09	
2.	Magnesium, Mg		0.02	675.0
3.	Aluminium, Al	mg/kg	269.972	53,900
4.	Potassium, K		0.12	11.10
5.	Calcium, Ca		0.09	293.0
6.	Chromium, Cr	mg/kg	2.18	14.40
7.	Manganese, Mn	mg/kg	19.36	3.99
8.	Iron, Fe	mg/kg	368.60	44500
9.	Nickel, Ni	mg/kg	1.81	28.90
10.	Copper, Cu	mg/kg	1.06	19.8
11.	Zinc, Zn	mg/kg	3.94	54.3
12.	Arsenic, As	mg/kg	ND<0.OOI	43.0
13.	Selenium, Se	mg/kg	ND<0.02	ND
14.	Silver, Ag	mg/kg	0.32	<0.5
15.	Cadmium, Cd	mg/kg	0.053	0.09
16.	Antimony, Sb	mg/kg	ND<0.OI	4.1E+02
17.	Lead, Pb	mg/kg	8.97	36.0
18.	Phosphorus, P		0.49	
19.	Mercury, Hg	mg/kg	ND<0.OOI	0.42
20.	pH		8.08	
21.	Total Nitrogen as Nitrogen		0.28	
21.	Total Organic Carbon		5.58	
23.	Moisture Content		6.43	

*Source: *Standard heavy metals limit based on the Department of Environment (DOE) Malaysia (2019)

This combined study helps to understand how the soil can be used in relation to environmental and agricultural activities (Eisenhauer et al., 2026). The macronutrient analysis of the soil shows that there are moderate amounts of vital nutrients in the soil, such as total nitrogen (0.28%), phosphorus (0.49%), and potassium (0.12%). Nonetheless, the potassium concentration is much lower than the DOE standard 11.10% indicating that it may require supplementation to facilitate plant growth. Particularly, a deficiency of potassium may impair the growth of plants and decrease their yields, which demonstrates the importance of specialised fertilisation. As Wang et al. (2026) define, potassium is an important macronutrient to plants, and one of the significant factors that inhibits plant growth is the lack of potassium. The lack of potassium in the soil results in a reduction in the size of the leaves, photosynthesis is reduced, and the leaf life is shortened. Severe potassium deficiency may lead to the turning of young leaves to pale yellow, and leaves become smaller and thinner.

The trace metal analysis indicates that the majority of the parameters are within DOE limits. As an example, magnesium (20 mg/kg), aluminium (269.972 mg/kg), and iron (368.60 mg/kg) are much lower than their corresponding limits of 675mg/kg, 53,900mg/kg, and 44,500mg/kg, making the soil safe for ecological use. Other trace metals, including nickel (1.81 mg/kg), copper (1.06 mg/kg), zinc (3.94 mg/kg), and silver (0.32 mg/kg), are well within permissible limits of 28.90 mg/kg, 19.8 mg/kg, 54.3 mg/kg, and 0.5 mg/kg, respectively. At the same time, lead (8.97 mg/kg) and antimony (< 0.01 mg/kg, non-detectable) also comply with DOE limits of 14.40 mg/kg and 410 mg/kg. In addition, toxic elements, including arsenic (< 0.001 mg/kg), selenium (< 0.02 mg/kg), cadmium (0.053 mg/kg), and mercury (< 0.001 mg/kg), are found in insignificant amounts, which cause minimal risks.

Nevertheless, manganese (19.36 mg/kg) is higher than the DOE value of 3.99 mg/kg, and it is concerning that it could cause environmental and agricultural effects. The high concentration of manganese could be explained by the natural sources in the sanctuary ecosystem, such as mineral-rich sediments and organic contributions, such as bird droppings. Furthermore, manganese toxicity could harm plant growth by making it harder for plants to absorb nutrients. Although this is a challenge, metagenomic data show that there are manganese-resistant microbial genes. This indicates that the microbial community of the soil may be essential in reducing the effect of this metal by converting it into less harmful forms. Besides industrial importance, Mn is an important element in the environment. It is one of the most important micronutrients, and almost all microbes require it. Manganese levels are increasing in the ecosystems across the world due to the high levels of industrial activity. Nevertheless, as they possess innate cellular processes that promote homeostasis throughout the ecosystem, bacterial strains play a vital role in the biogeochemical cycling of Mn (Ghosh & Das, 2018).

The TOC content of the soil is 5.58%, which indicates a healthy level of organic matter. Organic carbon plays a vital role in soil fertility, microbial activity and the cycling of nutrients. This high TOC indicates that the sanctuary environment, which is enhanced with organic inputs like bird activity and plant matter, is dynamic and resilient. Moreover, the metagenomic data demonstrated a diversity of microbes such as bacteria, archaea, Eukaryota and viruses as well as genes that are related to the cycling of nutrients, bioremediation and heavy metal resistance. This functional diversity indicates the influence of the physicochemical characteristics of the soil, including nutrient levels and heavy metal content, on the composition and activity of the microbial community (Dincă et al., 2022). The paper correlates physicochemical parameters with microbial functions to reveal the importance of soil chemistry in sustaining microbial diversity and ecological processes. This association is especially important when it comes to the UNISEL Bird Sanctuary, where the ecosystem, through its natural processes, provides inputs to the soil like organic matter through bird activity, which helps in the formation of the soil and its microbial community.

Quality Assessment and Metagenome Assembly Assessment

After filtering, a total of 28,583,695 reads with 62% GC content were obtained from whole-genome shotgun sequencing of the soil sample collected from UNISEL's Bird Sanctuary. The clean reads after the quality assessment were then used for metagenome assembly. It has successfully generated

2,850,843 contiguous sequences (contigs) numbers, with 439,071 for the largest contig, 451 bp for N50, and a total length of 1,245,901,409 bp, as summarised in Table 2.

Table 2
Metagenome assembly statistics

Contig no.	2,850,843
Largest contig	439,071
N50	451
Total length (bp)	1,245,901.409

Taxonomic Profiling

The high-throughput sequencing illustrates the diversity of the microbial community at the phylum level, as presented in the Krona plot in Figure 1. Figure 1(a) demonstrates the taxonomy profiling of microbial diversity in UNISEL's Bird Sanctuary using clean reads, while Figure 1(b) offers the taxonomy profiling based on the assembled data. In this study, clean reads were used directly for taxonomic profiling using Kraken 2 against the nt database. An average of 30.01% bacteria were observed through Kraken classification, with approximately 58.83% unclassified reads. Conversely, taxonomic profiling using assembled contigs showed an average of 34.83% bacteria, of which 51.68% were not classified. The high percentage of unclassified reads suggests the presence of possibly new microbial taxa and genetic functions, which are due to the underrepresentation of tropical soil microbiomes in reference databases and highlight the novelty and exploratory nature of this dataset (Ju & Zhang, 2015).

In the meantime, 8.38% of the reads (preassembly) were determined as Eukaryota, 0.21 as Archaea, 0.05 as viruses, and 2.52 as others. Taxonomic profiling of the post-assembly reads has identified 9.01% of the reads as Eukaryota, 0.25% as Archae, 0.05% as Viruses, and 4.18% as others. The difference between the percentages of bacteria in the clean reads and assembled contigs could be that the process of assembly enhanced the precision of taxonomic profiling. Conversely, assembly is reassembling longer genomic sequences out of smaller random derivatives. This may result in more precise taxonomic profiling, enabling longer and more informative sequences to be identified. It is important to compare the results of taxonomic profiling using clean reads and assembled contigs because they provide different information about the composition of a microbial community (Tran & Phan, 2020).

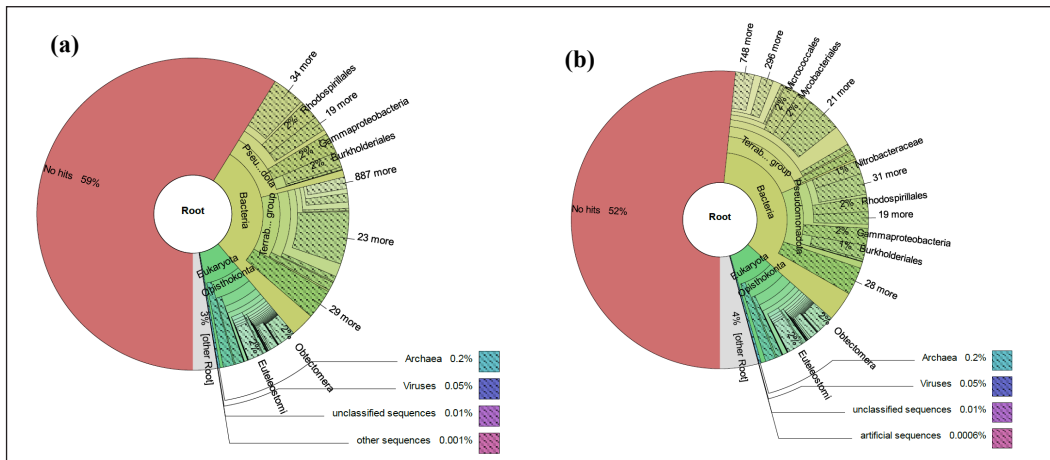


Figure 1. (a) Krona plot of profiled taxonomy of the clean reads (before assembly) against the Kraken 2 database; (b) Krona plot of profiled taxonomy of the assembled contigs against the Kraken 2 database

Kraken 2 is very reliable in identifying known taxa because it has an exact k-mer matching algorithm, which reduces false positives. Nevertheless, the sensitivity of this method is predetermined by the contents of world reference databases. The large percentage of unclassified sequences (51.7% - 58.8%) in this research is a direct result of the large underrepresentation of UNISEL bird sanctuary soil environments in genomic databases at the global and Malaysian levels, in particular. Since tropical freshwater sediments are still significantly underrepresented in metagenomic studies, a significant portion of our data corresponds to new microbial dark matter. This observation emphasises the significance of this preliminary baseline study, and by describing these uncharacterised communities for the first time in a Malaysian sanctuary, we are data mining a baseline that is needed to fill the gap between temperate-centric databases and tropical ecosystem reality.

Gene Prediction and Functional Annotation

The assembled contigs were then used for gene prediction using Prodigal (v2.6.3; -p meta -c -m), successfully identifying 3,282,381 genes with a total base of 348,958,57. The predicted genes were searched against the NCBI Reference Sequence (RefSeq) protein database, Swiss-Prot protein database, InterProScan, EggNOG, MEGARes, and KEGG database. Statistics of the annotations are summarised in Table 3 and Figures 2-5.

Figure 2 visually represents the InterPro annotation (Dimonaco et al., 2022; Hyatt et al., 2013) distribution within the soil metagenome dataset. The figure displays the various functional categories assigned to the predicted genes based on the InterPro databases. The highest contribution of the genes can be observed in the GO term database. This indicates that a significant portion of the predicted genes in the metagenome can be

assigned to broad functional categories related to biological processes, molecular functions, and cellular components. Meanwhile, moderate gene distribution was demonstrated in Superfamily, PANTHER, Gene3D, Pfam, MetaCyc, and Reactome databases, suggesting a substantial representation of genes involved in specific protein families, pathways, and reactions. Low gene distribution was annotated in SMART, CDD, ProSiteProfiles, ProSitePatterns, PRINTS, PIRSF, Coils, and MobiDBLite databases, where the lower representation of genes annotated by these databases could indicate that the functions associated with these specific domains and motifs are less prevalent in the microbial community in the studied area. Other than that, the InterPro annotation distribution provides a glimpse into the functional potential of the microbial community in the soil sample. It determines a diverse range of metabolic capabilities, emphasising functions related to broad GO categories and specific protein families and pathways represented by the moderately distributed databases.

Table 3
Functional annotation statistics against six different databases

Database	No. of Genes	% Genes
RefSeq	2,148,491	65.5
Swiss-Prot	1,176,492	35.8
InterProScan	2,070,549	63.1
eggNOG	2,200,678	67.0
MEGARes	2258	0.07
KEGG	1,024,042	31.2

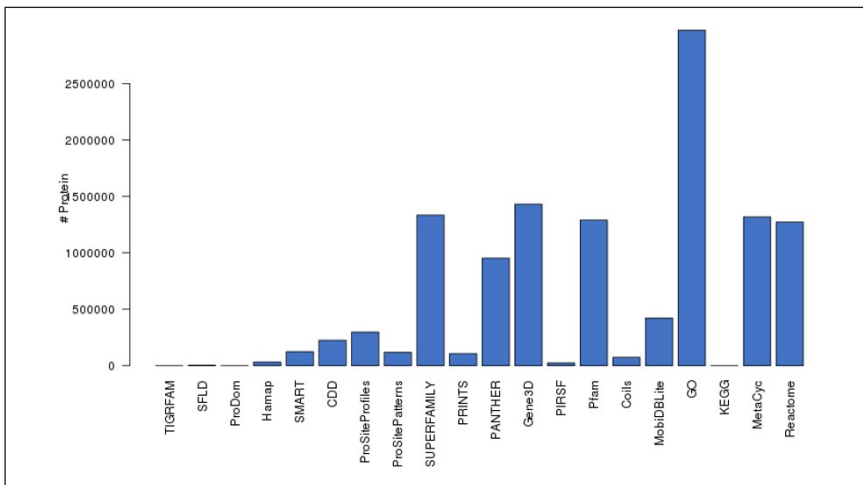


Figure 2. InterPro annotation showing the potential biological functions encoded by the microbial community in the UNISEL's Bird Sanctuary soil

The eggNOG analysis results illustrated in Figure 3 (a) suggest that bacteria are the primary contributors to the genetic diversity and functional potential of the microbial community in this environment. COG analysis, on the other hand, classifies predicted genes into broader functional categories to provide a general overview of the community's capabilities. Figure 3 (b) shows the COG classification analysis, which shows that the greatest percentage of genes is categorised as “Function Unknown.” This indicates that the microbial community in the soil sample contained many genes of unknown function. The large fraction of genes in the COG analysis that are categorised as Function unknown indicates a large gap in knowledge of the functional capabilities of the microbial community. This finding indicates that a significant fraction of the genes in the predicted metagenome of the soil sample does not have a distinct functional attribution according to the existing databases and knowledge.

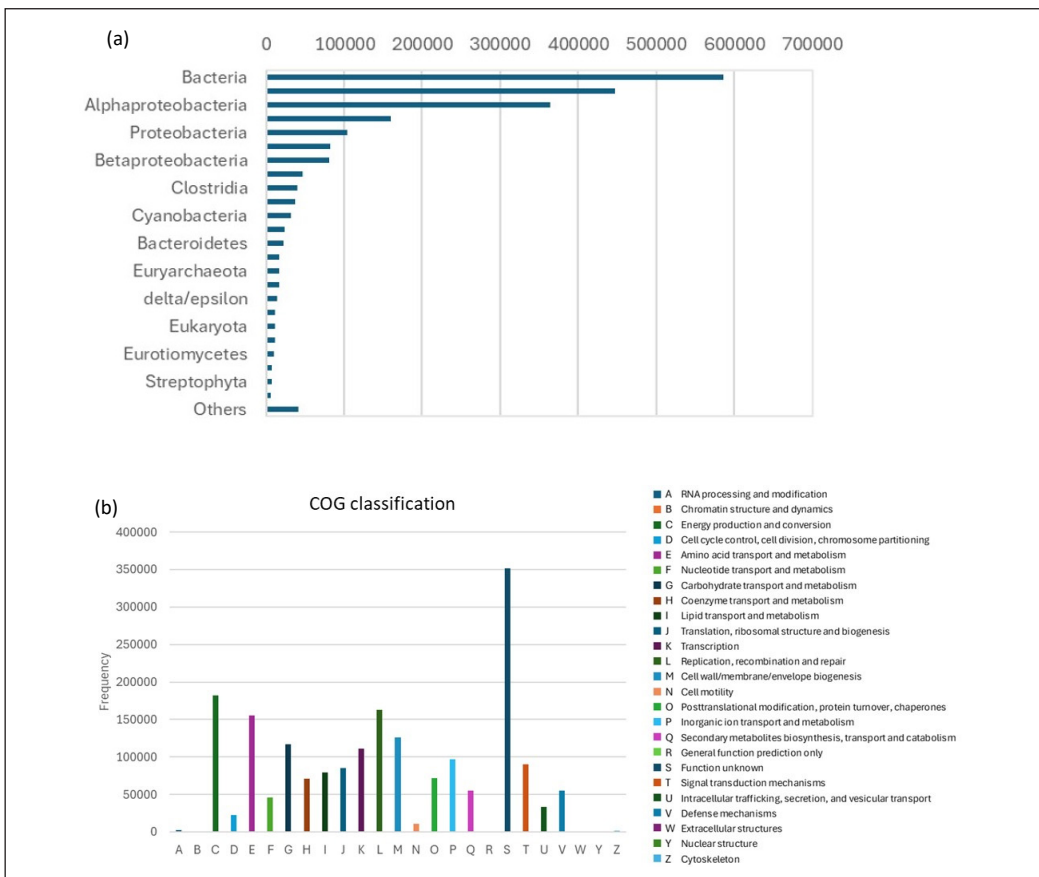


Figure 3. (a) The result from the EggNOG analysis shows the distribution of genes categorised under the bacteria group. (b) COG classification of the predicted genes presents the distribution of proteins based on their orthologous relationships

Likewise, Chen et al. (2022) stated that the role of soil microbial communities in ecosystem processes is essential in predicting how terrestrial ecosystems would act in response to climate change. The high proportion of unclassified sequences in the taxonomic analysis is consistent with the assignment of more than one-third of the predicted genes to 'function unknown' categories, which is in line with the high proportion of unclassified sequences found in the taxonomic analysis. Global functional databases are largely populated with sequences from temperate environments and clinical isolates, leaving tropical freshwater ecosystems like the Malaysian sanctuary soil significantly underrepresented.

The second most abundant gene set in the COG analysis, identified as an “Energy production and conversion” group, indicates that a significant percentage of the microbial community in the lake soil must be engaged in energy-related metabolic activities. This observation emphasises the importance of energy metabolism in the development of the community structure and functioning in this ecosystem (Chen et al., 2022). The third most common gene distribution is categorised as “Replication, recombination and repair”, and the fourth most common as “Amino acid transport and conversion”, which demonstrates some of the basics of the biology and activity of the lake soil microbial community. Essentially, all living organisms need the processes of replication, recombination, and repair to properly duplicate and preserve genetic material (Deng, 2023).

Further analysis of the predicted genes was conducted using MEGARes annotation to classify the presence of genes associated with drug resistance. Figure 4a shows that a significant percentage of genes are in the drug resistance category. This observation implies that the sanctuary soil sample contains a significant pool of genes that might potentially play a role in antimicrobial resistance. Particularly, the presence of these genes may be explained by the use of antibiotics or other antimicrobial agents in the environment. This brings into question the possibility of spreading resistance amongst microbes, which makes it difficult to treat infections caused by the microbes.

According to Table 1, physicochemical analysis indicates that the soil sample contains a high amount of iron. Thus, another annotation on MEGARes was conducted to study the genes around the metal resistance genes, and it showed the predominance of the genes around iron resistance (Figure 4b). The abundance of some metals in the soil sample and the metal resistance genes analysis with MEGARes are good evidence of the adaptation of the microbial community to this particular environment. Furthermore, the number of iron resistance genes is very high, which directly correlates with the high iron levels in the soil. This implies that the microbial community has developed ways of surviving and potentially uses iron, which could otherwise be harmful in high concentrations. According to Doster et al. (2020), MEGARes allow the detection of a wide range of metal resistance genes, which provides a broader understanding of microbial adaptation to the metal population.

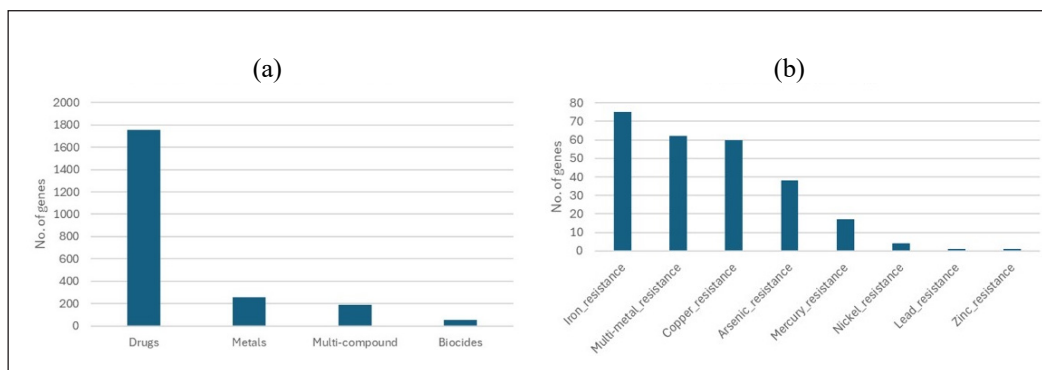


Figure 4. (a) MEGARes annotation of the assembled data of the soil sample, (b) Further MEGARes annotation on the metals-related genes of the soil sample

UNISEL Bird Sanctuary is found on a landscape that is historically related to tin mining operations. This can be seen in the physicochemical outcomes, which show high iron levels in the soil. Functional annotation of the metagenomic data also supported these conclusions by identifying genes linked to metal resistance-related genes, implying adaptation by the microbial adaptation to metal-enriched environments. The analysis of the metagenome data further indicated the most prevalent taxon, *Zygonium ericetorum*, reported to have extreme iron and aluminium stress tolerance (Herburger et al., 2016). The prevalence of these metal-tolerant microorganisms suggests an ecologically relevant role in mediating metal stress, leading to soil stability and allowing vegetation to establish in post-mining environments and could be investigated in future research.

MEGARes annotation also demonstrated that other genes are classified under multi-compound and biocide resistance categories. Multi-compound resistance genes confer resistance to multiple antimicrobial compounds, often with multiple mechanisms of action. These genes might encode efflux pumps that expel various compounds from the cell or enzymes that modify multiple drugs, rendering them ineffective (Maillard, 2018). Notably, biocide resistance genes enable microbes to withstand the effects of biocides, which are chemical agents used to control or kill harmful organisms (Maillard, 2018). These genes can provide resistance mechanisms like reducing biocide uptake, modifying the target site of the biocide, or detoxifying the biocide. According to Maillard (2018), efflux pumps are membrane protein which actively eliminate harmful compounds from cells, such as antibiotics and heavy metals, thus lowering intracellular concentrations and preventing harmful consequences.

Figure 5 shows the distribution of gene annotations according to the KEGG database. The KEGG database gathers manually curated pathway maps that reflect what we know about molecular interactions and reaction networks in biological systems (Kanehisa et al., 2023).

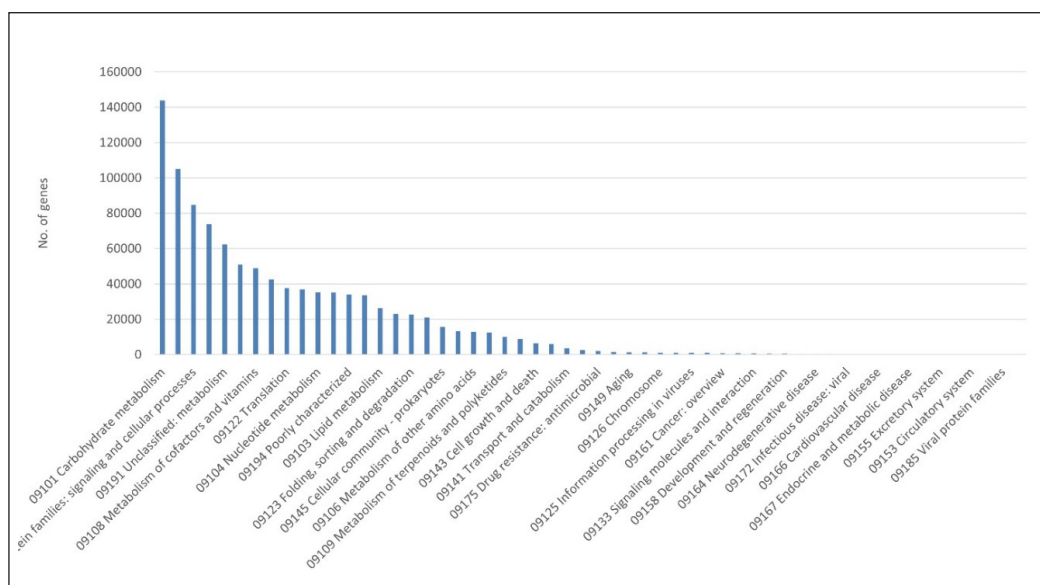


Figure 5. KEGG annotation based on the predicted genes identified in the studied soil sample

The figure is a visual representation of the percentage of genes present in the bird sanctuary soil sample related to different KEGG pathways. A major part of the genes is attributed to "Metabolism," which implies that the microbial community has a substantial role in metabolism. The other major groups are "Genetic Information Processing", "Environmental Information Processing", and "Cellular Processes", indicating the diverse functional ability of the microbial community in the bird sanctuary soil. The outcome of the KEGG analysis, as shown in Figure 5, illustrates the complexity and versatility of the microbial community in the way it performs a variety of biological processes. This involves metabolism, genetic information processing, environmental responses, and cellular functions, and this provides a variety of opportunities to conduct research in future studies.

Comparative Analysis

The UNISEL Bird Sanctuary soil metagenomic dataset was then compared to two publicly accessible NCBI datasets, which was done by a comparative analysis to provide context. These datasets were chosen as they are freshwater-affected soils and offer well-annotated shotgun metagenomic data and can be comparatively analysed taxonomically and functionally. The bird sanctuary soil sample was compared with two public datasets, which were received in BioProject PRJNA889691 and reported by Ren et al. (2017). Henceforth, the soil sample at the Bird Sanctuary of UNISEL will be referred to as S3. Figure 6 shows a phylum-level taxonomy of the three sites (S3, SRR21870203, and SRR21870204),

demonstrating similar and different microbial communities. The common phyla were Actinomycetota, Pseudomonadota, Bacillota, and Bacteroidota, which were all found with varying levels of abundance. The Actinomycetota dominated at S3. Furthermore, the microbial communities at SRR21870203 and SRR21870204 were dominated by Pseudomonadota. This implies that Pseudomonadota proliferation is supported by similar environmental factors or substrates at these locations. Although in varying amounts, Bacillota and Bacteroidota were constantly found at all sites.

Figure 7 shows the COG (Galperin et al., 2025) classification analysis of the three sites (S3, SRR21870203, and SRR21870204), indicating common and distinct functional profiles of microbial communities in the three sites. The three sites are highly represented in the major functional categories, including E (Amino acid transport and metabolism), G (Carbohydrate transport and metabolism), C (Energy production and conversion) and J (Translation, ribosomal structure, and biogenesis), which suggests their core functions in the maintenance of microbial activity. The categories such as R (General function prediction only) and S (Function unknown) are also prevalent, indicating the presence of genes that are either uncharacterised or predictive, which is a common characteristic of metagenomic studies.

S3 has a reduced total abundance of functional genes, indicating a less metabolically diverse community. This is especially noticeable in the categories like N (Cell motility) and T (Signal transduction mechanisms), which show lower representation, meaning fewer cellular interactions and environmental responses. By contrast, SRR21870203 and SRR21870204 show more functional diversity, with increased gene abundances in functional categories like L (Replication, recombination, and repair),

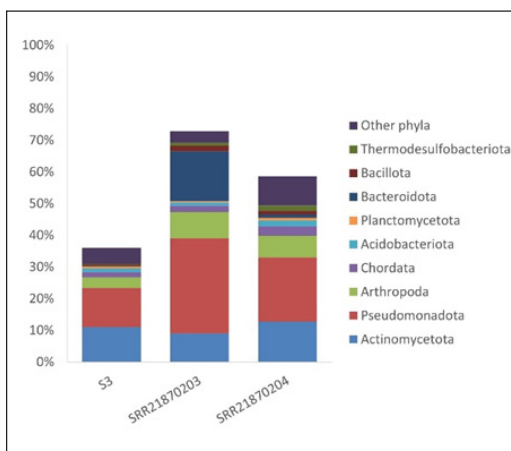


Figure 6. Comparison of the taxonomic profile of phyla in the three samples (S3, SRR21870203 and SRR21870204)

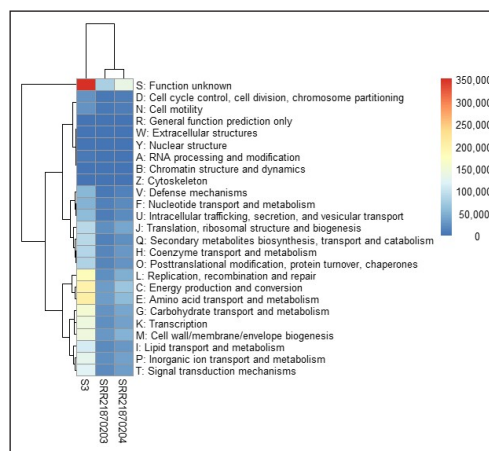


Figure 7. Comparison of COG classification in the three samples (S3, SRR21870203 and SRR21870204)

M (Cell wall/membrane/envelope biogenesis), and V (Defence mechanisms), reflecting more active cellular processes and enhanced adaptability. SRR21870204 also indicates increased representation in D (Cell cycle control, cell division, and chromosome partitioning) and Z (Cytoskeleton), indicating structural and division-related adaptations in this microbial community. Nevertheless, the fact that many uncharacterised genes in S (Function unknown) are highly abundant in S3 contributes to the possibility of new discoveries and the necessity of conducting additional studies to reveal the unknown functions of microbes (Cecil et al., 2018; Price et al., 2018).

Figure 8 shows the outcome of comparing the metal resistance genes in the three samples (S3, SRR21870203, and SRR21870204) using MEGARes annotation. This indicates variations in the prevalence of resistance to different metals due to site-specific environmental factors and possible contamination. Moreover, arsenic, copper, iron, nickel, zinc and lead resistance genes were found in all three samples, suggesting that all three were exposed to these metals, perhaps through contamination or natural geochemical conditions. According to Gillieatt and Coleman (2024), heavy metals are anthropogenic and natural in different sources and highly influential because of their presence in the environment. S3 had the greatest number of multi-metal resistance genes, indicating it was exposed to a combination of metals, with iron resistance genes being the most abundant. Essentially, the difference in the metal resistance genes highlights the importance of the local geochemical environment in shaping microbial populations.

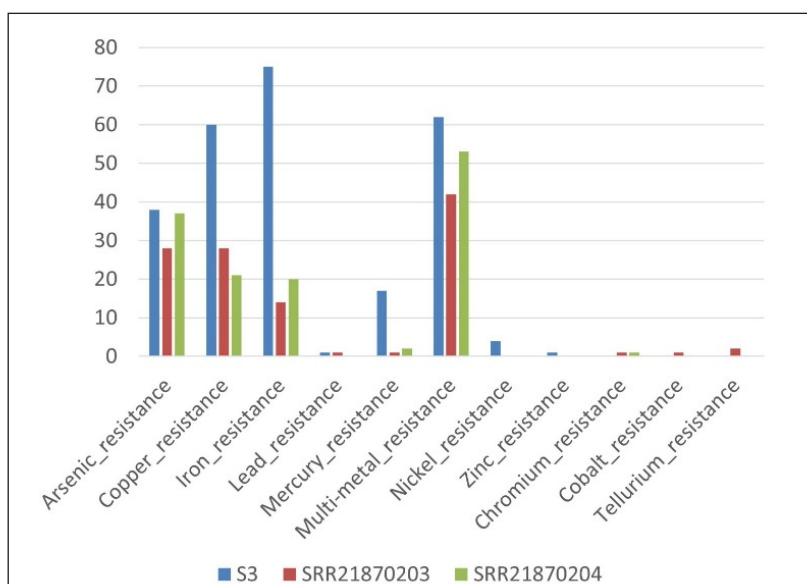


Figure 8. Comparative abundance of metal resistance genes

The PCA plot in Figure 9 represents the beta diversity analysis of microbial community composition across three samples: S3, SRR21870204, and SRR21870203. Beta diversity quantifies the difference in the microbial community composition among samples, and the PCA plot visualises these differences across two principal components (V1 and V2). The sampling distribution pattern shows that there are different microbial compositions (Mikha et al., 2024). Sample S3 is located separately from SRR21870204 and SRR21870203, suggesting a separate microbial community. Also, SRR21870204 is more distant from S3 and SRR21870203, implying that it possesses the most unique microbial composition. Remarkably, S3 is closer to SRR21870203 than to SRR21870204, which means that S3 is more similar to SRR21870203. The differences in microbial diversity observed might be explained by multiple ecological, environmental, or experimental factors. These can be differences in habitat conditions (e.g., soil or water environment), time variability in sampling, or possible technical differences in processing or sequencing procedure (Jiao et al., 2022).

In this research, CCA was used to determine how the microbial taxonomy is related to soil parameters in the three samples (Figure 10) in order to determine the effects of the soil characteristics on different microbial phyla (Bai et al., 2017; Chen et al., 2022; Zhou et al., 2023). The position of the microbial taxa (marked in red) and the three samples with these axes shows the correlation of the microbial communities with particular variables of the soil. SRR21870204 is located far away from SRR21870203 and S3, showing a different microbial composition under unique soil conditions. In contrast, SRR21870203 and S3 are more similar, which implies more similar microbial communities, probably due to similar soil factors. The position of microbial taxa close to a particular sample suggests

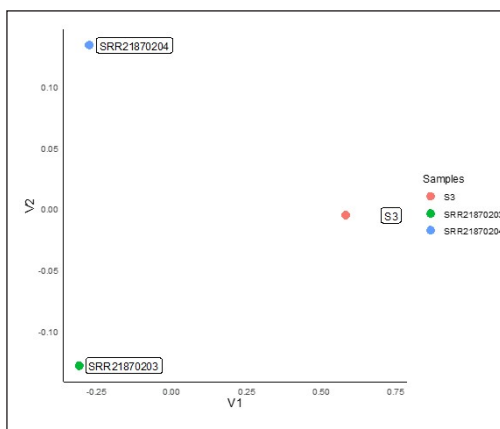


Figure 9. PCA plot (beta-diversity) of taxonomic phyla

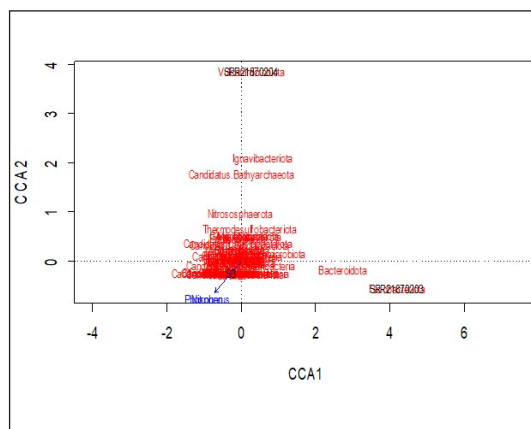


Figure 10. CCA analysis of taxonomy against soil phyla

a close relationship with the microbial community of the sample (He et al., 2023; Sari et al., 2024). In particular, Ignavibacteriota and *Candidatus* Bathyarchaiota are closer to SRR21870204, implying that these phyla are more common in that sample. The blue arrow marked “Phthhokens” is a variable of soil that affects the distribution of microbes. The direction and length of the arrow show how strong this effect is, and the taxa or samples that are closer to the arrow are more affected by this aspect.

The comparative metagenomics indicated similarities and differences between the soil sample and other public datasets, particularly in taxonomic composition, COG classification, and the abundance of metal resistance genes. The increased concentration of copper, iron and mercury resistance genes in the bird sanctuary soil sample indicates potential adaptations to certain environmental conditions. Generally, the metagenomic analysis offers a comprehensive description of the microbial population within the soil, including its diversity and functional attributes.

CONCLUSION

The research gives a detailed baseline characterisation of soil microbial communities in the UNISEL Bird Sanctuary through a combination of shotgun metagenomics and physicochemical methodology. The findings provide a record of the native microbiome of the indigenous soil of a minimally disturbed tropical ecosystem, which will provide useful reference material to future ecological, environmental and conservation research, and also to conservation natural heritage records. The structure of the microbial community and functional profiles were closely correlated with the physicochemical characteristics of the soil, especially mineral composition and indicated the adaptation of microorganisms to geochemical and historical conditions of the site. The abundance of metal-tolerant taxa and related functional genes emphasises the ecological importance of soil microorganisms in ecosystem stability and vegetation establishment in a post-mining landscape.

The metagenomic dataset showed a significant percentage of unclassified reads, which probably explained the underrepresentation of tropical soil microbiomes in the available reference databases, rather than issues with data quality. The observation indicates the novelty of the UNISEL Bird Sanctuary soil microbiome and implies the presence of previously unexplored microbial taxa and genetic functions. These unclassified reads highlight the value of such a dataset as a baseline reference data and demonstrate a large potential in the future of discovery with improved annotation resources and increased sequencing depth. However, additional studies are necessary to investigate the functions and interactions of the uncharacterised microbial components and adaptive mechanisms in this ecosystem, including additional sequencing and other highly advanced analysis methods.

ACKNOWLEDGEMENT

Our deepest gratitude goes to Universiti Selangor (UNISEL) for trusting and funding this research project under the Geran Penyelidikan Dalaman Universiti Selangor (BESTARI) 2023 (GPB/UNISEL-23/ST/01).

REFERENCES

- Bai, R., Wang, J., Deng, Y., He, J., Feng, K., & Zhang, L. (2017). Microbial community and functional structure significantly varied among distinct types of paddy soils but responded differently along gradients of soil depth layers. *Frontiers in Microbiology*, 8, Article 945. <https://doi.org/10.3389/fmicb.2017.00945>
- Cecil, J. H., Garcia, D. C., Giannone, R. J., & Michener, J. K. (2018). Rapid, parallel identification of catabolism pathways of lignin-derived aromatic compounds in *Novosphingobium aromaticivorans*. *Applied and Environmental Microbiology*, 84(22), Article e01185-18. <https://doi.org/10.1128/AEM.01185-18>
- Chen, H., Ma, K., Lu, C., Fu, Q., Qiu, Y., Zhao, J., Huang, Y., Yang, Y., Schadt, C. W., & Chen, H. (2022). Functional redundancy in soil microbial community based on metagenomics across the globe. *Frontiers in Microbiology*, 13, Article 878978. <https://doi.org/10.3389/fmicb.2022.878978>
- Crawford, K. M., Dice, C. G., & Clark, G. S. (2026). Why, when, and how microbes can benefit ecological restorations: current approaches and future directions. *New Phytologist*. <https://doi.org/10.1111/nph.70995>
- Deng, S. (2023). The origin of genetic and metabolic systems: Evolutionary structural insights. *Heliyon*, 9(3). <https://doi.org/10.1016/j.heliyon.2023.e14466>
- Department of Environment (DOE) Malaysia (2019). *Contaminated land management and control guidelines no. 1: Malaysian recommended site screening levels for contaminated land*. Ministry of Environment and Water. https://www.doe.gov.my/wp-content/uploads/2021/07/Contaminated-Land-Management-and-Control-Guidelines-No-1_Malaysian-Recommended-Site-Screening-Levels-for-Contaminated-Land.pdf
- Dimonaco N. J., Aubrey W., Kenobi K., Clare A., Creevey C. J. (2022). No one tool to rule them all: Prokaryotic gene prediction tool annotations are highly dependent on the organism of study. *Bioinformatics*, 38(5), 1198-1207. <https://doi.org/10.1093/bioinformatics/btab827>
- Dincă, L. C., Grenni, P., Onet, C., & Onet, A. (2022). Fertilisation and soil microbial community: A review. *Applied Sciences*, 12(3), Article 1198. <https://doi.org/10.3390/app12031198>
- Doster, E., Lakin, S. M., Dean, C. J., Wolfe, C., Young, J. G., Boucher, C., Belk, K. E., Noyes, N. R., & Morley, P. S. (2020). MEGARes 2.0: a database for classification of antimicrobial drug, biocide and metal resistance determinants in metagenomic sequence data. *Nucleic acids research*, 48(D1), D561-D569. <https://doi.org/10.1093/nar/gkz1010>
- Eisenhauer, N., Sünemann, M., Pollierer, M. M., Sun, X., Bardgett, R. D., Bartkowski, B., Delgado-Baquerizo, M., Dirilgen, T., Guerra, C. A., Mathieu, J., Niklaus, P. A., Ristok, C., Seeber, J., Steinwandter, M., Stewart, J. Heijden, M. V. D., Putten, W. V. D. and Potapov, A. (2026). Soil biodiversity effects on ecosystems. *Nature Reviews Biodiversity*, 1-16. <https://www.nature.com/articles/s44358-025-00123-z>
- Fulke, A. B., Mahajan, M. S., & Gothwal, S. K. (2026). Insights on the applicability of metagenomics in marine environments: ecological and biotechnological perspectives. *Marine Biology*, 173(1), 1. <https://doi.org/10.1007/s00227-025-04735-z>

- Galperin, M. Y., Vera Alvarez, R., Karamycheva, S., Makarova, K. S., Wolf, Y. I., Landsman, D., & Koonin, E. V. (2025). COG database update 2024. *Nucleic Acids Research*, *53*(D1), D356-D363. <https://doi.org/10.1093/nar/gkae983>
- Ghosh, S., & Das, A. P. (2018). Metagenomic insights into the microbial diversity in manganese-contaminated mine tailings and their role in biogeochemical cycling of manganese. *Scientific Reports*, *8*(1), Article 8257. <https://doi.org/10.1038/s41598-018-26311-w>
- Gillicatt, B. F., & Coleman, N. V. (2024). Unravelling the mechanisms of antibiotic and heavy metal resistance co-selection in environmental bacteria. *FEMS Microbiology Reviews*, *48*(4), Article fuac017. <https://doi.org/10.1093/femsre/fuac017>
- He, Z., Yuan, C., Chen, P., Rong, Z., Peng, T., Farooq, T. H., Wang, G., Yan, W., & Wang, J. (2023). Soil microbial community composition and diversity analysis under different land use patterns in Taojia River Basin. *Forests*, *14*(5), Article 1004. <https://doi.org/10.3390/f14051004>
- Herburger, K., Remias, D., & Holzinger, A. (2016). The green alga *Zygonium ericetorum* (Zygnematophyceae, Charophyta) shows high iron and aluminium tolerance: Protection mechanisms and photosynthetic performance. *FEMS Microbiology Ecology*, *92*(8), Article fiw103. <https://doi.org/10.1093/femsec/fiw103>
- Husain, R., Vikram, N., Pandey, S., Yadav, G., Bose, S. K., Ahamad, A., Shami, M., Kumar, K. G., Khan, N.A., Zuhaib, M., & Hussain, T (2022). Microorganisms: An eco-friendly tools for the waste management and environmental safety. In S. Arora, A. Kumar, S. Ogita, & Y. Y. Yau (Eds.), *Biotechnological innovations for environmental bioremediation* (pp. 949-981). Springer. https://doi.org/10.1007/978-981-16-9001-3_36
- Hyatt, D., Chen, G. L., LoCascio, P. F., Land, M. L., Larimer, F. W., & Hauser, L. J. (2010). Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC bioinformatics*, *11*(1), 119. <https://doi.org/10.1186/1471-2105-11-119>
- Jiao, S., Chu, H., Zhang, B., Wei, X., Chen, W., & Wei, G. (2022). Linking soil fungi to bacterial community assembly in arid ecosystems. *iMeta*, *1*(1), Article e2. <https://doi.org/10.1002/imt2.2>
- Ju, F., & Zhang, T. (2015). Experimental design and bioinformatics analysis for the application of metagenomics in environmental sciences and biotechnology. *Environmental Science & Technology*, *49*(21), 12628-12640. <https://doi.org/10.1021/acs.est.5b03719>
- Kanehisa, M., Furumichi, M., Sato, Y., Kawashima, M., & Ishiguro-Watanabe, M. (2023). KEGG for taxonomy-based analysis of pathways and genomes. *Nucleic Acids Research*, *51*(D1), D587-D592. <https://doi.org/10.1093/nar/gkac963>
- Kumar, R. (2026). Isolation and cultivation of novel microorganisms for drug prospecting. In *Bioinformatics, AI, and Machine Learning in Microbial Drug Development* (pp. 3-22). Academic Press. <https://doi.org/10.1016/B978-0-443-33032-2.00012-5>
- Li, L., Ye, J., Yu, M., Jiang, J., Guo, X., Yu, W., & Rong, K. (2025). Dynamic changes in the avian gut microbiome in response to diverse lifestyles. *Ibis*, *167*(2), 331-344. <https://doi.org/10.1111/ibi.13388>
- Maillard, J. Y. (2018). Resistance of bacteria to biocides. *Microbiology Spectrum*, *6*(2), Article arba-0006-2017. <https://doi.org/10.1128/microbiolspec.ARBA-0006-2017>

- McLaren, M. R., Hershey, O. S., Machtinger, A. N., Rice, D. P., Simas, A. M., Friedman, C. R., Gratalo, D., Philipson, C. W., & Bradshaw, W. J. (2026). Metagenomic sequencing of composite airplane wastewater for surveillance of emerging viruses. *medRxiv*. <https://doi.org/10.64898/2026.01.29.26343714>
- Mikha, M. M., Green, T. R., Untiedt, T. J., & Hergret, G. W. (2024). Land management affects soil structural stability: Multi-index principal component analyses of treatment interactions. *Soil and Tillage Research*, 235, Article 105890. <https://doi.org/10.1016/j.still.2023.105890>
- Price, M., Wetmore, K., Waters, R., Callaghan, M., Ray, J., Liu, H., Kuehl, J., Melnyk, R., Lamson, J., Suh, Y., Carlson, H., Esquivel, Z., Sadeeshkumar, H., Chakraborty, R., Zane, G., Rubin, B., Wall, J., Visel, A., Visel, A., Bristow, J., Blow, M., Arkin, A., Arkin, A., Deutschbauer, A., & Deutschbauer, A. (2018). Mutant phenotypes for thousands of bacterial genes of unknown function. *Nature*, 557(7706), 503-509. <https://doi.org/10.1038/s41586-018-0124-0>
- Ren, Z., Wang, F., Qu, X., Elser, J. J., Liu, Y., & Chu, L. (2017). Taxonomic and functional differences between microbial communities in Qinghai Lake and its input streams. *Frontiers in Microbiology*, 8, Article 2319. <https://doi.org/10.3389/fmicb.2017.02319>
- Salam, L. B., & Obayori, O. S. (2026). Shotgun metagenomics reveals the complete genetic potential for lindane biodegradation in a tropical lentic pond sediment. *International Biodeterioration & Biodegradation*, 206, Article 106206. <https://doi.org/10.1016/j.ibiod.2025.106206>
- Sari, O. F., Bader-El-Den, M., Ince, V., & Arabikhan, F. (2024). *Machine learning approach into bacterial relationship: Exploring 16S rRNA metabarcoding with association rule mining* [Paper presentation]. 2024 IEEE 12th International Conference on Intelligent Systems (IS), Varna, Bulgaria. <https://doi.org/10.1109/IS61756.2024.10705245>
- Tran, Q., & Phan, V. (2020). Assembling reads improves taxonomic classification of species. *Genes*, 11(8), Article 946. <https://doi.org/10.3390/genes11080946>
- Wang, Z., Yao, C., Yang, Y., Segbo, S., Xu, X., Lin, X., Zhou, P., Gao, F., Ni, Z., Shi, T., & Gao, Z. (2026). Effects of foliar potassium fertiliser on photosynthetic capacity and expression of potassium and sugar transporters in peach (*Prunus persica*). *Horticulturae*, 12(3), Article 388. <https://doi.org/10.3390/horticulturae12030388>
- Xi, D., Zhu, F., Zhang, Z., Zhou, S., & Zhang, J. (2026). Forest type shapes soil microbial carbon metabolism: A metagenomic study of subtropical forests on Lushan Mountain. *Microorganisms*, 14(1), Article 220. <https://doi.org/10.3390/microorganisms14010220>
- Zainun, M. Y., & Simarani, K. (2018). Metagenomics profiling for assessing microbial diversity in both active and closed landfills. *Science of the Total Environment*, 616, 269-278. <https://doi.org/10.1016/j.scitotenv.2017.10.266>